

Increasing TCP's Initial Window

draft-hkchu-tcpm-initcwn-01.txt

Nandita Dukkkipati

Yuchung Cheng

Jerry Chu

Matt Mathis

{nanditad, ycheng, hkchu, mattmathis}@google.com

30 July, 2010

78th IETF, Maastricht

Overview of prior results for IW10

- Our proposal: Increase TCP's IW to 10 MSS
- IW10 improves average TCP latency by ~10%
- Large scale data-center experiments demonstrate latency improves across network and traffic properties:
 - varying network bandwidths, flow RTTs, bandwidth-delay products, HTTP response sizes, mobile networks
 - small overall increase in retransmission rate (~0.5%), with most from multiple connections
- Prior work:
 - <http://www.ietf.org/proceedings/10mar/slides/tcpm-4.pdf>
 - <http://ccr.sigcomm.org/online/?q=node/621>

New contributions and the questions addressed

- A framework for running experiments with different IWs in the same data-center
- Primary concern from IETF-77: how does IW10 perform on highly multiplexed links such as in Africa and South America?
- What is the impact on latency due to losses in the initial window?
- Evaluated the impact of different IWs [3, 10, 16] on latency and retransmission rate
 - Reinforced the prior experiment results with IW10
- Testbed experiments for IW study in controlled environment
 - Preliminary results on fairness

Improved methodology for experiments

Previous methodology:

- Change IW for entire data-center every week
 - Less apples-to-apples: changes in server software and user base
 - Takes weeks to collect data

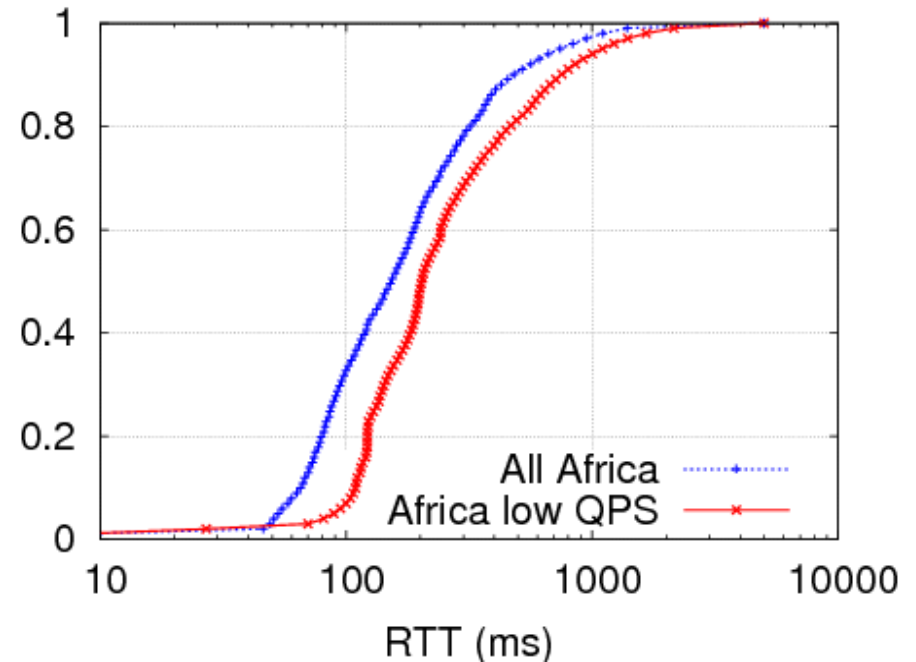
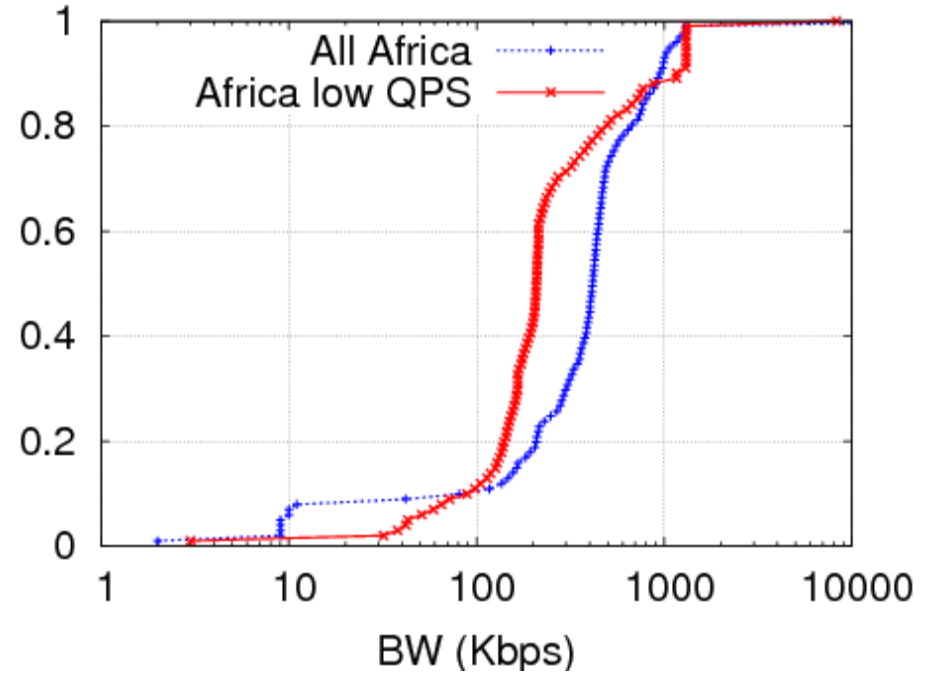
New methodology:

- Serve different IWs based on IP address in one data-center simultaneously for weeks
 - Same IW for connections from the same IP/browser
 - More apples-to-apples: free from binary/config changes

Analysis of IW10 on Africa traffic



Experiment for 1 week in
June 2010



Impact of IW10 on Africa traffic

Web Search latency (ms) and retransmission rate %

All of Africa

Percentile	Avg.	50	75	90	99
IW=10	988.4	503	795	1467	5042
IW=3	1123.9	538	878	1710	5923
Impr.	135.5	35	83	243	881
% Impr.	12%	6.5%	9.5%	14.2%	14.9%

	Retrans. %
IW=10	3.77%
IW=3	3.35%
Increase	0.42

Africa with low QPS

Percentile	Avg.	50	75	90	99
IW=10	1870.5	733	1363	3146	11579
IW=3	2340.7	857	1773	4110	14414
Impr.	470.2	124	410	964	2835
% Impr.	20.1%	14.5%	23.1%	23.5%	19.7%

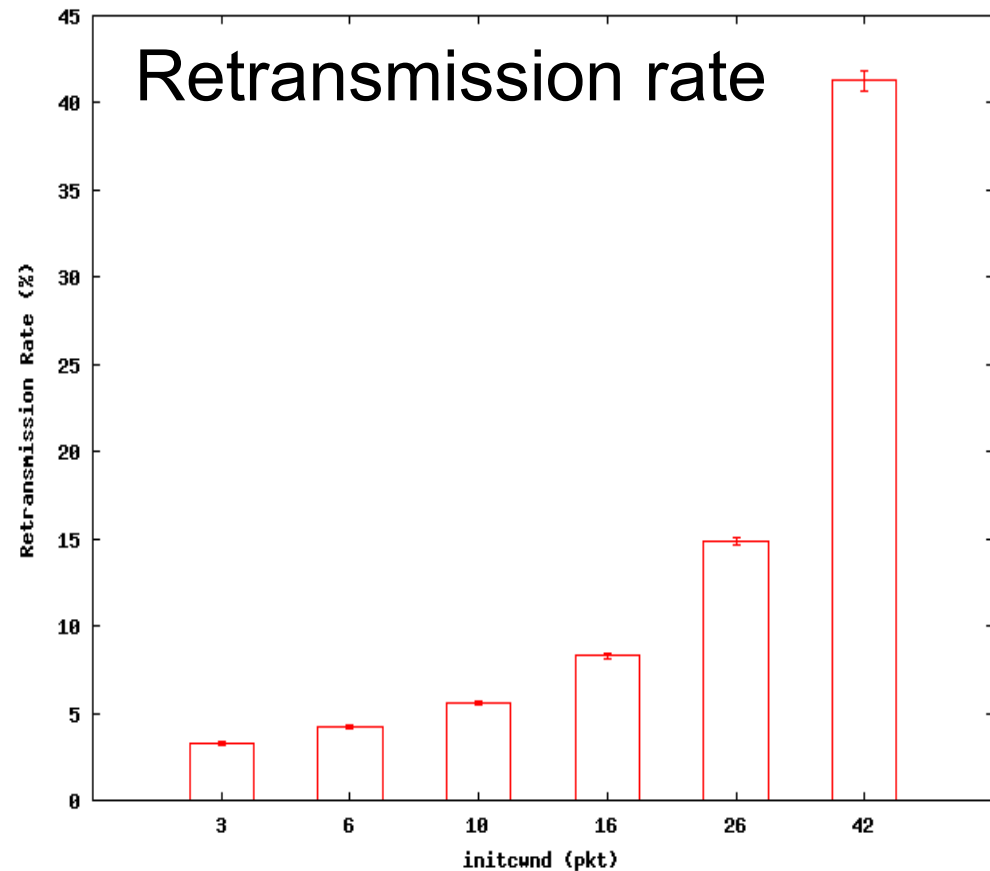
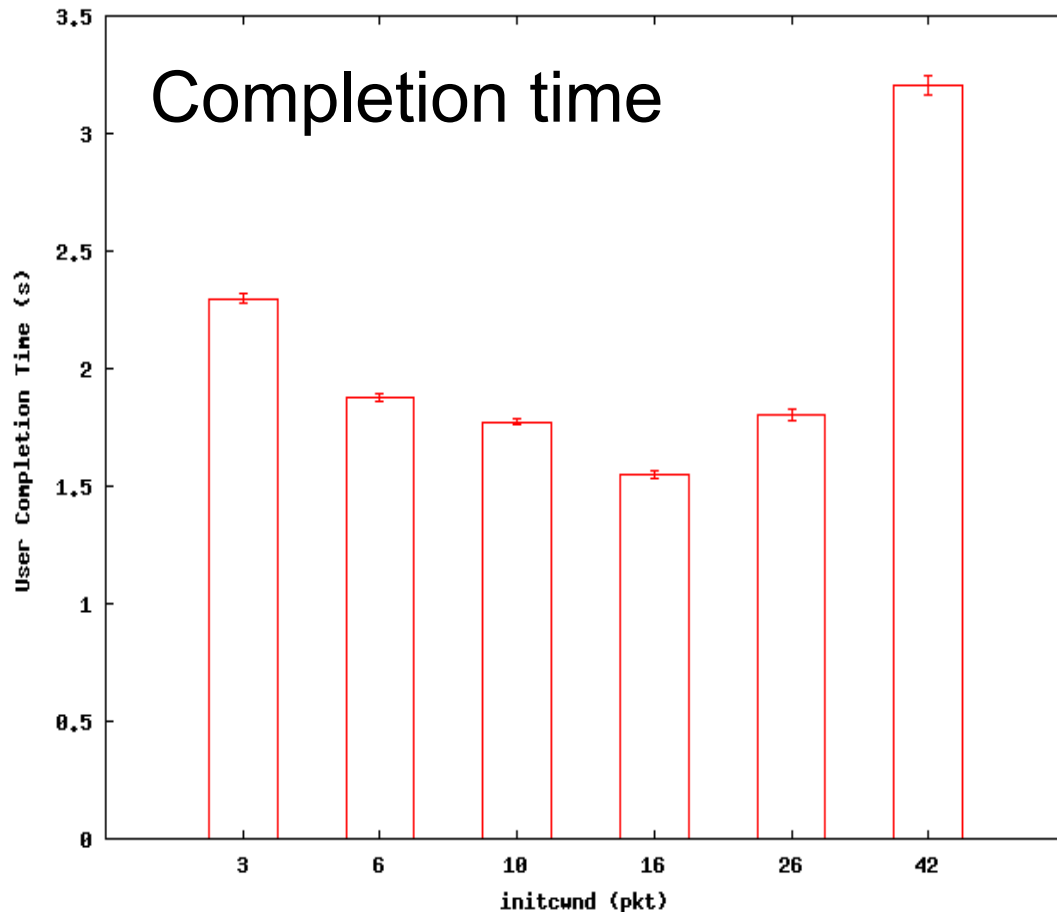
	Retrans. %
IW=10	6.71%
IW=3	5.83%
Increase	0.87

Why does latency improve in Africa?

- Large network round-trip time
- Larger IW helps faster recovery of packet losses
- Experiments on testbed demonstrate latency improves in spite of increased packet losses

Why does latency improve in Africa?

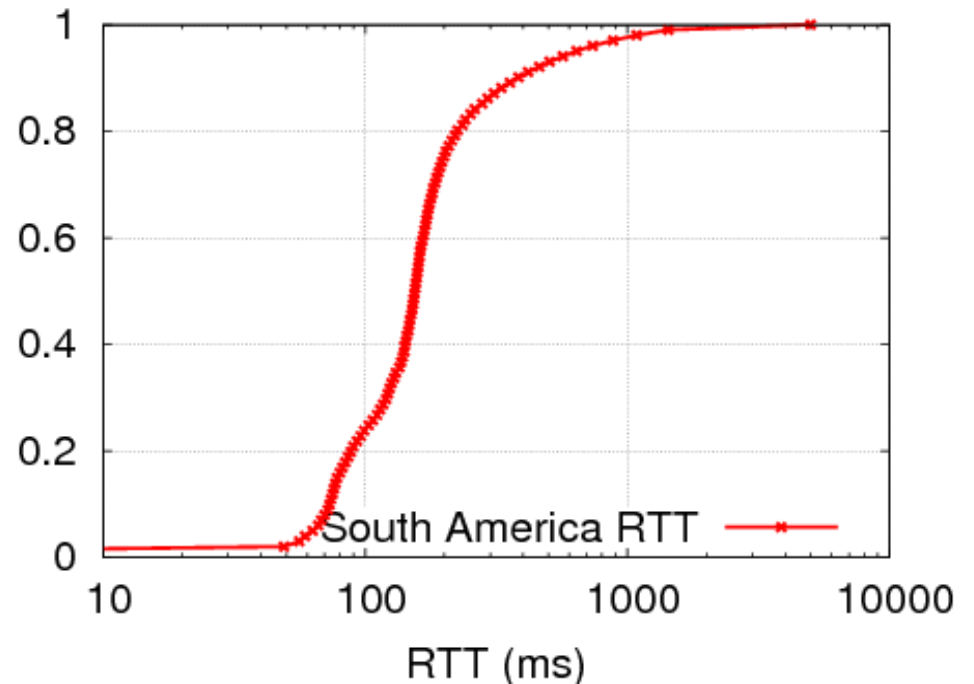
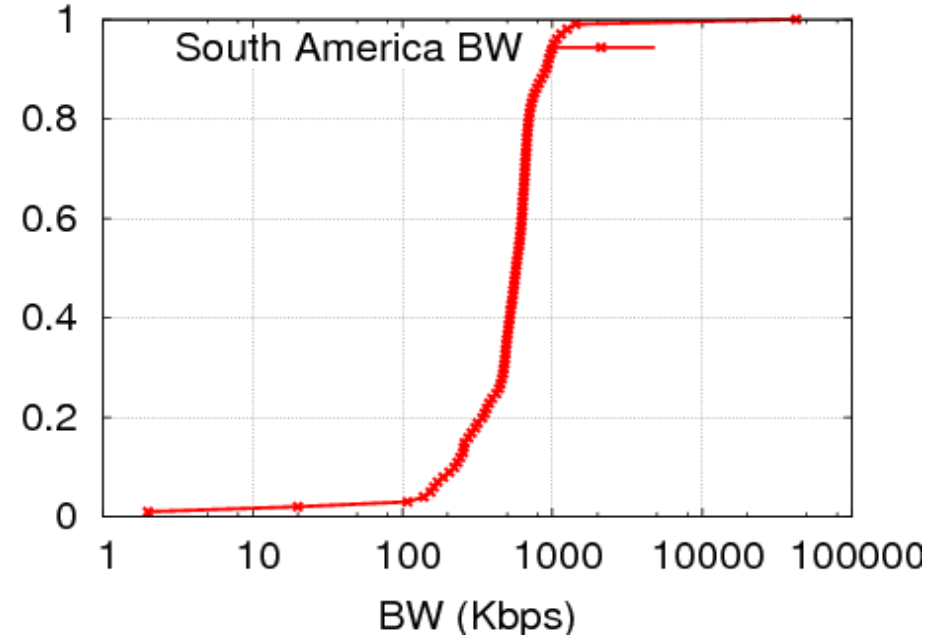
- Testbed experiment: 20Mbps, RTT 300ms, BDP buffer, offered load 0.95, 50KB response size
- Motivating example: Makerere University, Uganda



Analysis of IW10 on South America traffic



Experiment for 1 week in June 2010



South America latency with IW10

Web Search latency (ms) and retransmission rate %

All of South America

Percentile	Avg.	50	75	90	99
IW=10	919.7	540	1245	1970	4059
IW=3	1018.6	578	1399	2257	4620
Impr.	98.8	38	154	287	561
% Impr.	9.7%	6.6%	11%	12.7%	12.1%

	Retrans. %
IW=10	2.81%
IW=3	2.3%
Increase	0.51

Latency improvement across services in South America

- Latency improves across a variety of services
- Services with multiple connections experience:
 - Least latency benefits
 - Most increase in retransmission rate

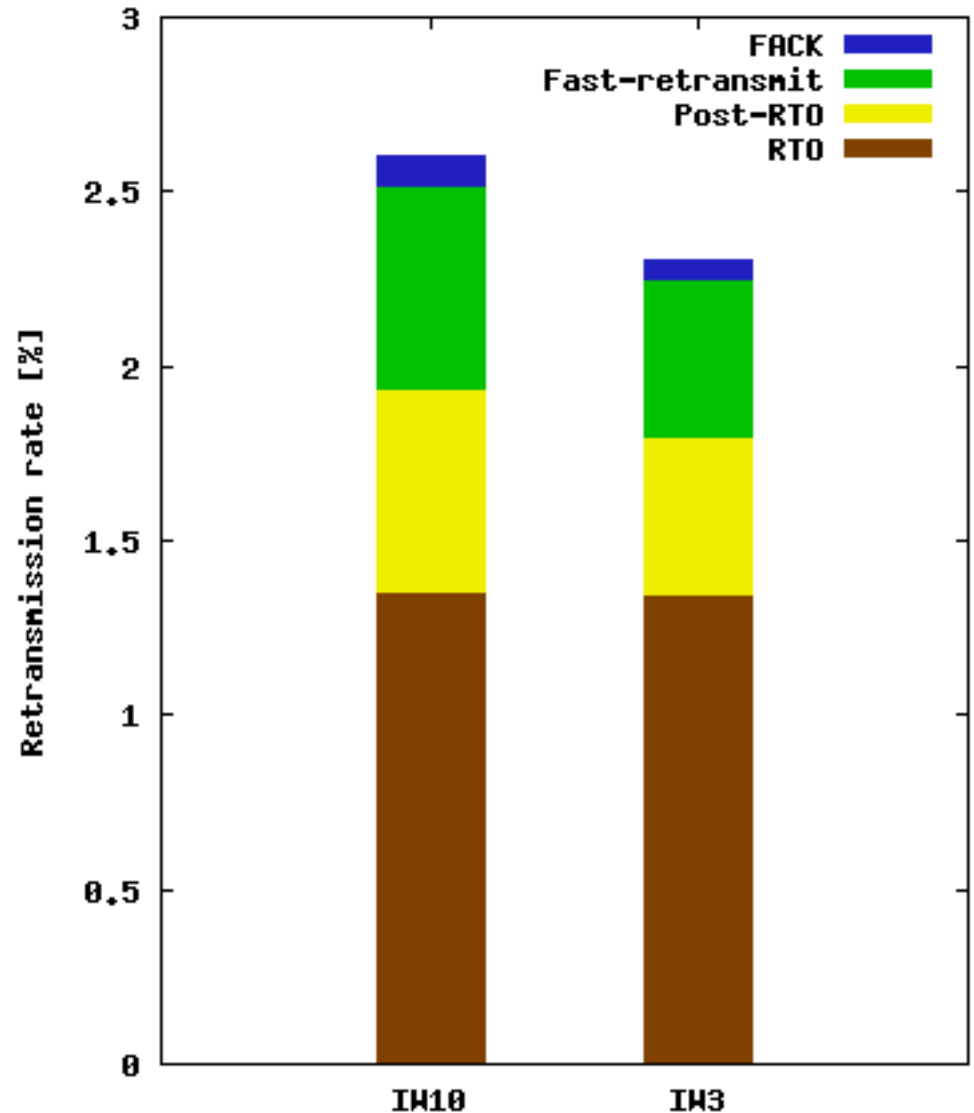
Percentile	iGoogle	News	Blogger Photos (multiple connections)	Maps (multiple connections)
10	30 [10%]	4 [2.5%]	2 [1.1%]	6 [3.8%]
50	198 [26%]	45 [9.9%]	98 [12.7%]	12 [3.2%]
90	430 [16%]	336 [15%]	251 [4.5%]	37 [2.6%]
99	986 [9.7%]	1827 [19%]	691 [2.9%]	134 [2.9%]
Delta in Retrans %	0.52	0.35	2.93	1.28

entry: latency improvement (ms) [% improvement]

Retransmission of IW3 vs IW10

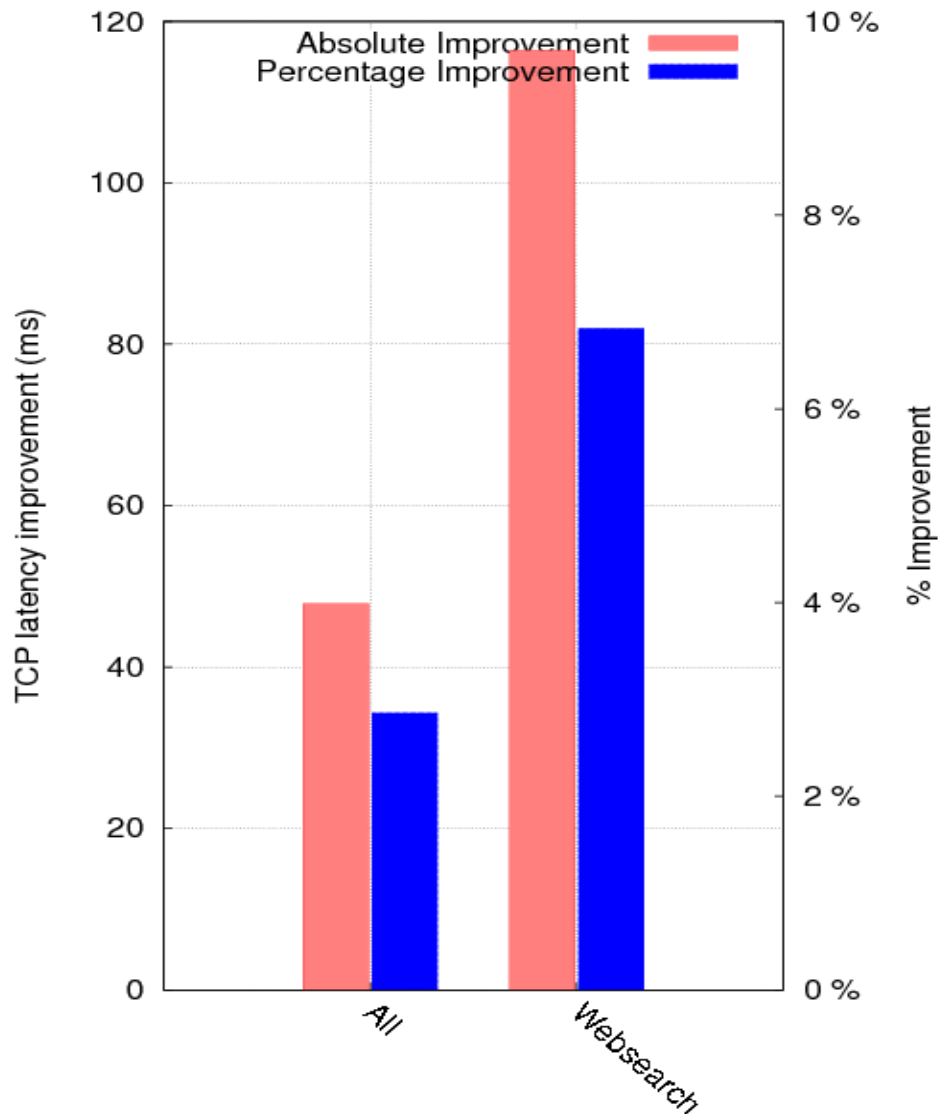
IW10 has no significant increase in timeouts, but has more

- fast-retransmit
- post-RTO retransmits



Impact of latency under packet losses

Latency of traffic with retransmissions > 0 improves with IW10 as compared to IW3



% traffic with retransmit > 0

	IW3	IW10
All	6.6%	6.8%
Web Search	6.11%	6.57%

Experiments with higher IWs

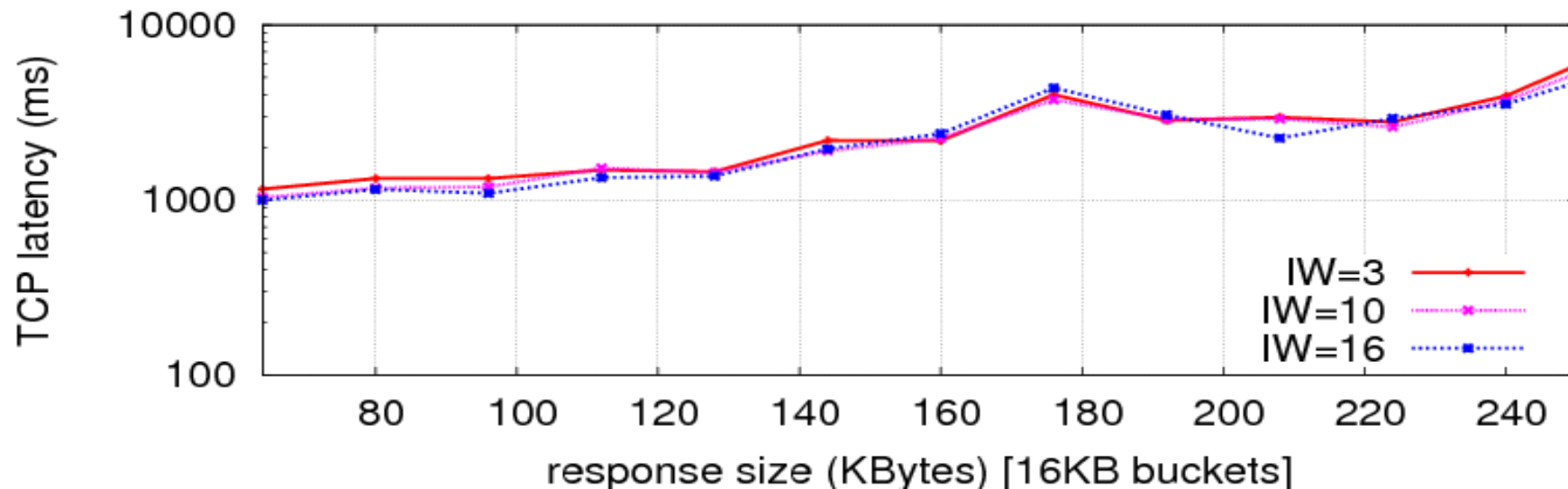
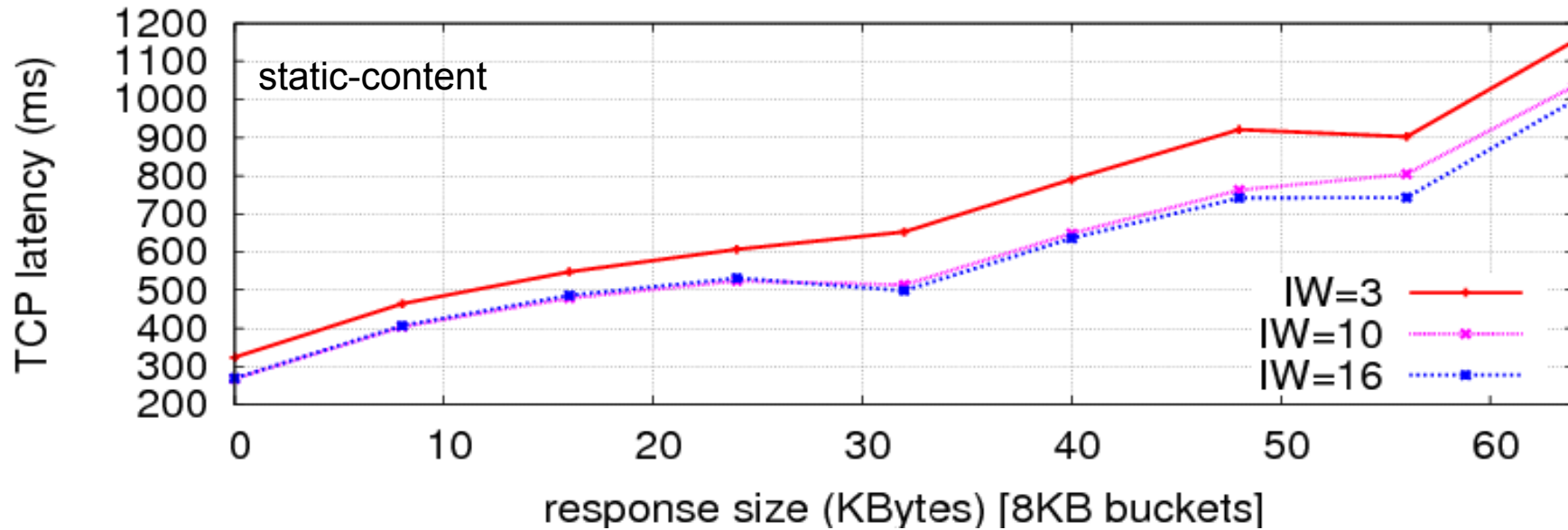
Does higher IW show better latency? What is the "sweet spot" of IW?

Client IP based IW Experiments:

- DC 1
 - 20% in US east coast (RTT < 100ms)
 - 80% in south America (RTT > 100ms)
- DC 2
 - 97% in Europe (RTT < 100ms)

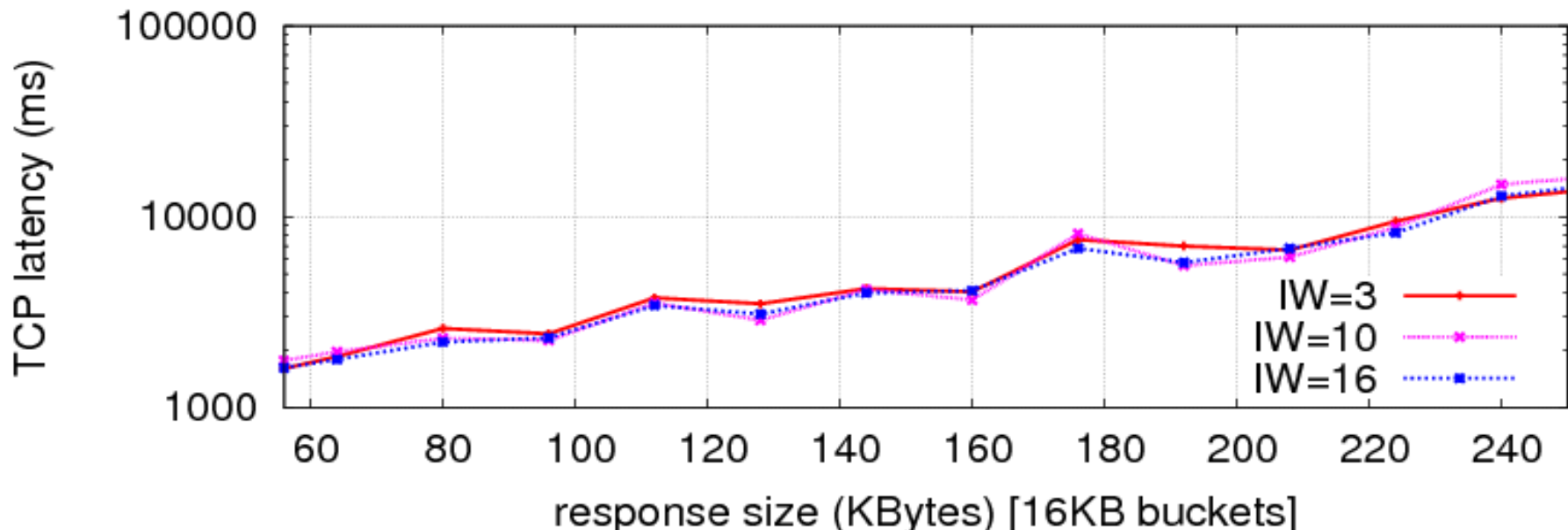
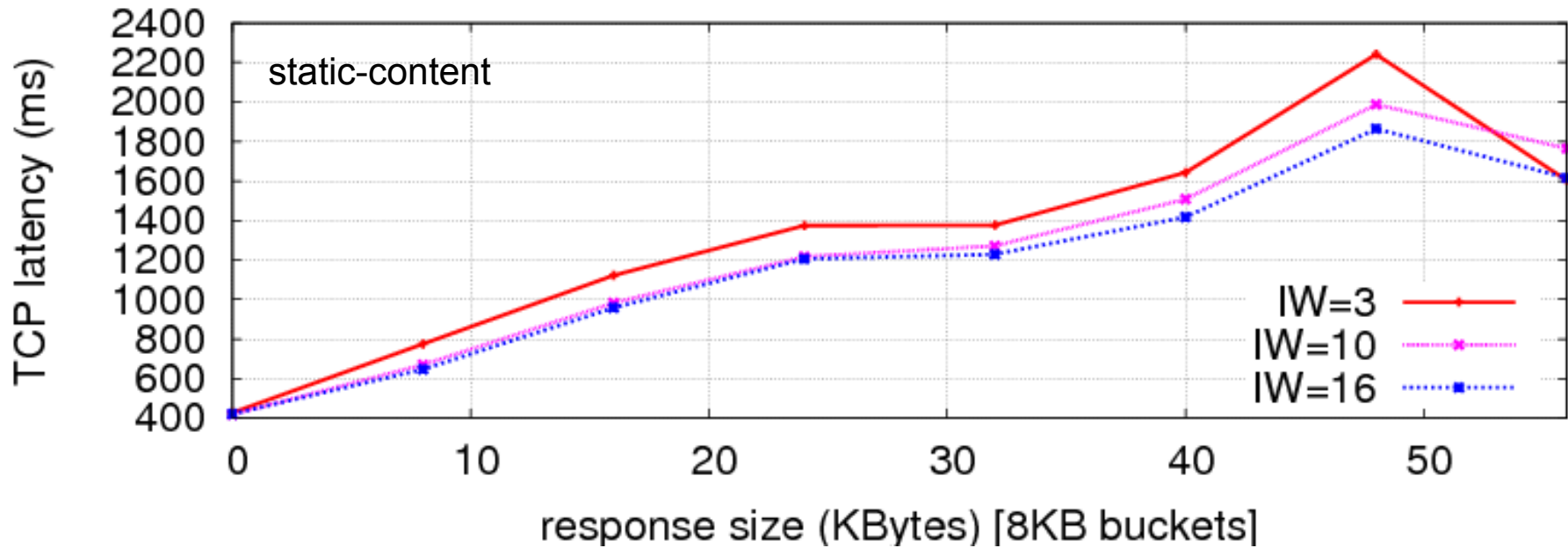
Comparison of IW = 3, 10, 16 (DC 1)

Small improvement for larger IWs (>10); mostly for mid-size flows



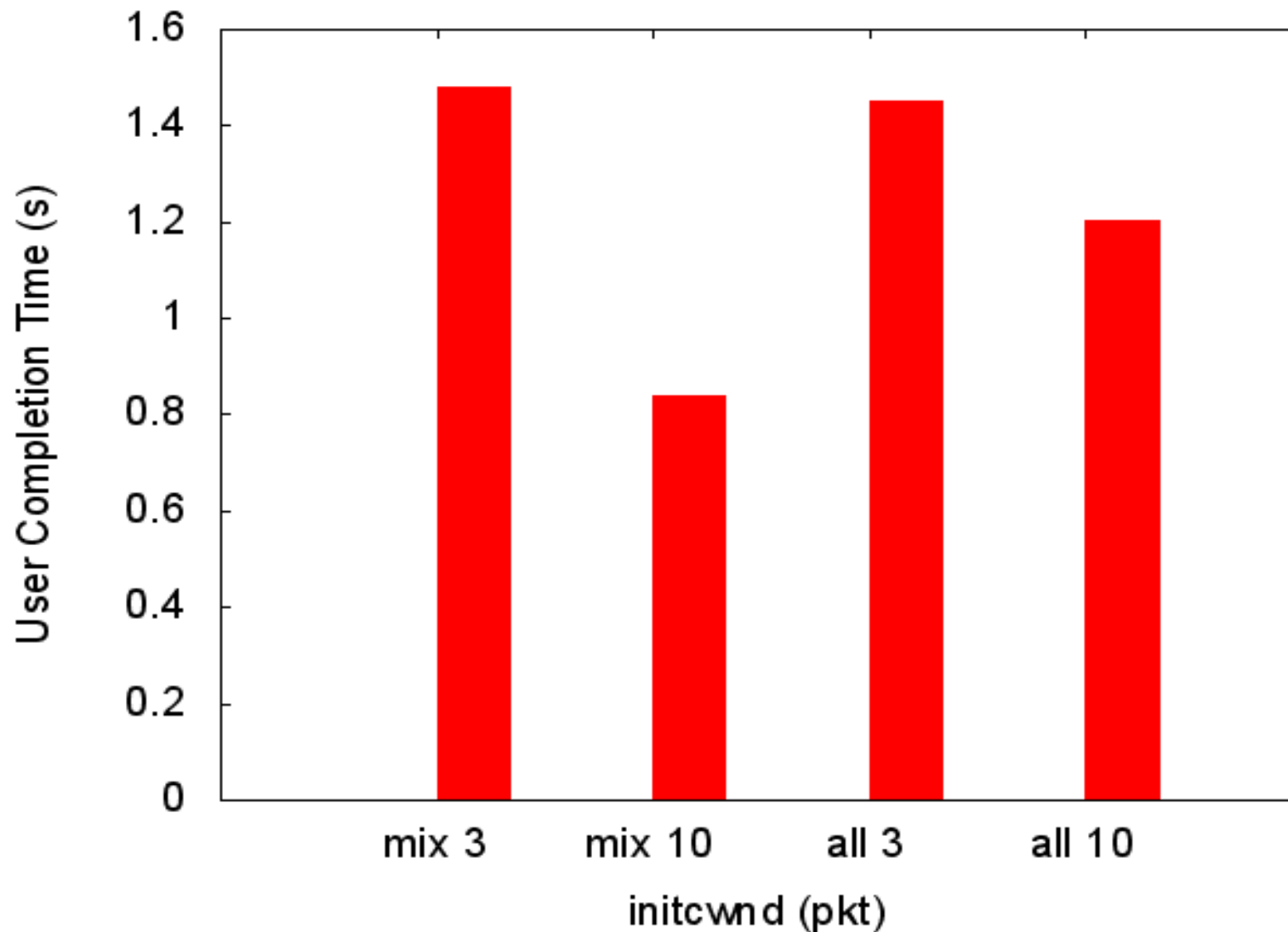
Comparison of IW = 3, 10, 16 (DC 2)

Small improvement for larger IWs (>10); mostly for mid-size flows



Fairness between IW10 and IW3 flows

Testbed experiment: 20Mbps, 300ms, BDP buffer, load 0.95, 15KB response size, mix of IW3 and IW10 traffic



Conclusion

- Take away summary
 - IW10 improves latency even in Africa and South America
 - IW10 helps in quicker recovery from packet losses
 - A higher retransmission rate does not necessarily translate to a longer Web transfer latency
 - IW16 shows a small latency improvement over IW10
- Next steps
 - Adoption of IW10 proposal as TCPM WG item
 - Ongoing work: fairness between IW3 and IW10 in the transition phase
 - For any pending issues with IW10, join us in solving the problems!

Steps to configure IW on Linux

Changing TCP IW on Linux (kernel version \geq 2.6.30)

On your server, do

\$ ip route show

select the outgoing route then do

\$ ip route change default via <gateway> dev eth0 initcwnd <iw>

If the server process explicitly set SNDBUF, then SNDBUF value \geq IW*MSS. Otherwise increase the initial socket buffer if IW*MSS > /proc/sys/net/ipv4/tcp_wmem[1]

\$ cat /proc/sys/net/ipv4/tcp_wmem

4096 16384 4194304

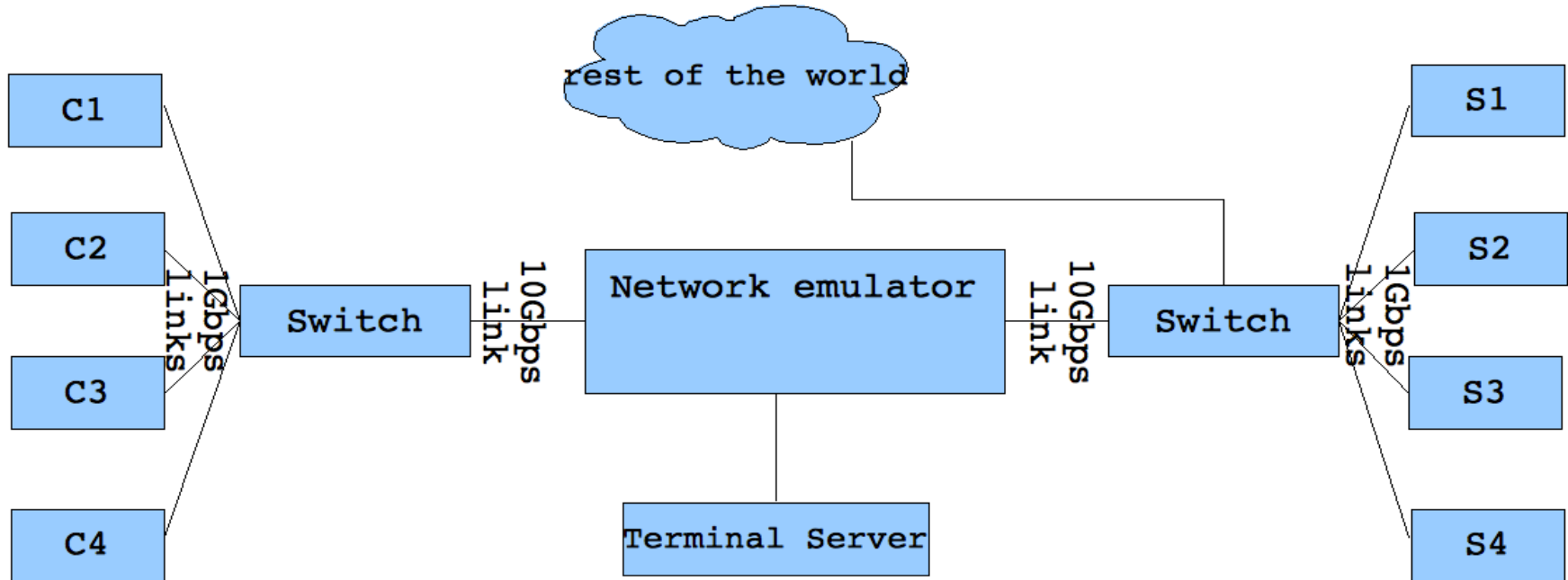
\$ echo '4096 IW*MSS 4194304' > /proc/sys/net/ipv4/tcp_wmem

\$ must restart server process to use new tcp_wmem[1]

Acknowledgements

- We acknowledge the following people at Google for their contributions towards the large scale IW experiments:
 - Ethan Solomita
 - Elliott Karpilovsky
 - John Reese
 - Yaogong Wang
 - Roberto Peon
 - Arvind Jain

Testbed topology



Why does latency improve in Africa?

- (Contd') Tesbed experiment results

